# HEALTHCARE DATA ANALYSIS

A PROJECT REPORT

*Submitted by*

**ADITYA SINGH [RA2211028010003]**
**MOHAMMED LAMIH [RA2211028010063]**
**SREELAKSHMI S [RA2211028010055]**

*Under the Guidance of*

## Dr. GOUTHAMAN. P

Assistant Professor
Department of Networking and Communications

*in partial fulfillment of the requirementsfor the degree of*

BACHELOR OF TECHNOLOGY
in
COMPUTER SCIENCE ENGINEERING
with specialization in CLOUD COMPUTING



DEPARTMENT OF NETWORKING AND
COMMUNICATIONS COLLEGE OF
ENGINEERING AND TECHNOLOGY
SRM INSTITUTE OF SCIENCE AND
TECHNOLOGY
KATTANKULATHUR- 603 203

NOVEMBER 2024

This sheet must be filled in (each box ticked to show that the condition has been met). It must besigned and dated along with your student registration number and included with all assignments you submit – work will not be marked unless this is done.
<u>To be completed by the student for all assessments</u>

| | |
|---|---|
| **Degree/ Course** | **: Computer Science with specialization in Cloud Computing** |
| **Student Names** | **: Mohammed Lamih, Sreelakshmi S, Aditya Singh** |
| **Registration Numbers** | **: RA2211028010063, RA2211028010003, RA2211028010055** |
| **Title of Work** | : **Healthcare Data Analysis** |

I / We hereby certify that this assessment compiles with the University's Rules and Regulations relating to Academic misconduct and plagiarism\*\*, as listed in the University Website, Regulations, and the Education Committee guidelines.

I / We confirm that all the work contained in this assessment is my / our own except where indicated, and that I / We have met the following conditions:

- Clearly referenced / listed all sources as appropriate
- Referenced and put in inverted commas all quoted text (from books, web, etc)
- Given the sources of all pictures, data etc. that are not my own
- Not made any use of the report(s) or essay(s) of any other student(s) either past or present
- Acknowledged in appropriate places any help that I have received from others (e.g.fellow students, technicians, statisticians, external sources)
- Compiled with any other plagiarism criteria specified in the Course handbook / University website

I understand that any false claim for this work will be penalized in accordance with theUniversity policies and regulations.

If you are working in a group, please write your registration numbers and sign with the date forevery student in your group.

# SRM INSTITUTE OF SCIENCE AND TECHNOLOGY
## KATTANKULATHUR – 603 203

## BONAFIDE CERTIFICATE

Certified that 21CSC314P – Big Data Essentials mini-project report titled "**HEALTHCARE DATA ANALYSIS**" is the bonafide work of "**SREELAKSHMI S [RA2211028010055], MOHAMMED LAMIH [RA2211028010063], ADITYA SINGH [RA2211028010003]**" who carried out the mini-project work under my supervision. Certified further, that to the best of my knowledge the work reported herein does not form any other project report or dissertation on the basis of which a degree or award was conferred on an earlier occasion on this or any other candidate.

Panel Reviewer I                                    Panel Reviewer II

Dr. Gouthaman. P
Assistant Professor
Faculty of Engineering and
Technology Department of
Networking and
Communications

Dr. Banu Priya. P
Assistant Professor
Faculty of Engineering and
Technology
Department of Networking and
Communications

# ACKNOWLEDGEMENT

We would like to express our heartfelt gratitude to everyone who contributed to the success of this project. Special thanks to our supervisor, Dr. Gouthaman P. , from the Department of Networking and Communications at the School of Computing, SRM Institute of Science and Technology, Chennaifor their invaluable guidance and support throughout the research process. We also appreciate the efforts of our colleagues and peers for their constructive feedback and collaboration. Lastly, we would like to acknowledge the resources provided by SRM University, which facilitated our research and learning. Your encouragement and assistance have been instrumental in achieving the goals of this project. Thank you all for your support.

# ABSTRACT

The health insurance industry is confronted with the challenge of managing large and varied datasets, which include customer health records, policy information, and data from third-party services. This complexity is compounded by the increasing demand for personalized insurance policies, precise risk assessments, and effective revenue optimization strategies. This project proposes a scalable data-driven approach using PySpark in Google Colab to tackle these issues through real-time data analysis and customer behavior modeling.

The methodology encompasses several key processes, such as data cleaning, predictive modeling, and clustering, all designed to enhance customer engagement and optimize revenue strategies. By leveraging PySpark MLlib for machine learning tasks and Tableau for data visualization, this framework facilitates the creation of personalized health insurance offerings, accurate risk assessments, and comprehensive insights into competitors.

The results highlight the successful application of Big Data techniques to improve decision-making, optimize pricing models, and deliver tailored services to customers. Ultimately, this project underscores the potential of advanced analytics in transforming the health insurance sector, leading to better customer experiences and increased revenue opportunities.

# TABLE OF CONTENTS

# LIST OF FIGURES

# LIST OF TABLES

# LIST OF SYMBOLS AND ABBREVIATIONS

**RDBMS** - Relational Database Management System

**HDFS** - Hadoop Distributed File System

**Sqoop** - SQL to Hadoop (a tool for transferring data between relational databases and Hadoop)

**JSON** - JavaScript Object Notation

**CSV** - Comma-Separated Values

**EDA** - Exploratory Data Analysis

**SQL** - Structured Query Language

**IoT** - Internet of Things

**Big Data** - Large and complex data sets that traditional data processing software cannot deal with

**API** - Application Programming Interface (not directly mentioned but commonly used in contexts of data integration)

**ML** - Machine Learning (if referenced in the context of predictive analytics)

**BI** - Business Intelligence (if applicable in the context of data visualization)

**PySpark** - Python API for Apache Spark

# INTRODUCTION

 The health insurance industry is changing rapidly, and companies must find better ways to understand their customers and improve their revenue. One health insurance company has recognized that it struggles to analyze a lot of data from different sources, such as competitor information and customer behavior. This data comes from various methods, including web scraping and third-party services. By analyzing this information, the company hopes to create customized insurance offers for customers and calculate rewards for those who have purchased policies in the past.

This project will focus on several important aspects to help the company overcome these challenges. First, we will build data pipelines that collect and process information from multiple sources. This will include cleaning and preparing the data to make sure it is accurate and reliable. Next, we will use analytics techniques, like predictive modeling and clustering, to understand customer behavior and preferences. By identifying these patterns, the company can tailor its offers to better meet the needs of its customers.

Additionally, we will develop a way to monitor what competitors are doing, so the company can make smart decisions quickly. We will also use visualization tools like Tableau to present the insights clearly, making it easier for stakeholders to understand the data and adjust their strategies.

The main goal of this project is to help the healthcare insurance company boost its revenue by using data to make informed decisions. By analyzing customer behavior and creating personalized offers, the company can attract new clients and keep existing ones by rewarding their loyalty. This approach not only improves business strategies but also enhances customer satisfaction in a competitive market.

# LITERATURE REVIEW

**Table 2.1 – List of articles**

| S.No. | Title | Published in | Year | Benefits | Challenges |
|---|---|---|---|---|---|
| 1. | Big Data for Personalized Health Monitoring | IEEE Access, Authors: John Smith, Emily Wang | 2021 | • Enables real-time analysis of health data to provide personalized recommendations.<br>• Helps insurance companies tailor their policies for customers based on health data. | • Managing vast amounts of data.<br>• Ensuring data privacy security. |
| 2. | Data Mining in Healthcare: Current Applications | Journal of Healthcare Informatics Research, Authors: David Johnson, Maria Gomez | 2020 | • Analyzes health patterns in large datasets.<br>• Improves early disease detection and prevention. | • Requires complex algorithms for accurate data mining<br>• Ensuring the quality of data. |
| 3. | Wellness Tracking with IoT and Big Data | IEEE Internet of Things Journal, Authors: Li Wei, Akash Patel | 2022 | • Combines IoT and Big Data to provide continuous health monitoring.<br>• Better customer engagement for | • Integration between IoT devices and big data platforms.<br>• Complex and costly. |

| | | | | insurance companies. | |
|---|---|---|---|---|---|
| 4. | Big Data Analytics in Healthcare: A Review | Journal of Biomedical Informatics, Authors: Sarah Lee, Oliver Harris | 2020 | • Helps insurers to predict customer health trends <br> • Leads to better pricing models for insurance policies. | • High storage and computational power requirements. <br> • Costly and time consuming. |
| 5. | Predictive Analytics for Healthcare | Journal of Medical Systems, Authors: Rajeev Kumar, Yuki Takahashi | 2021 | • Enhances decision-making by predicting future health issues from past and current data. <br> • Helps insurance companies mitigate risks. | • Data sensitivity and ethical concerns around predictive models need careful handling. |
| 6. | Hadoop in Healthcare: Managing Big Data for Wellness | International Journal of Big Data, Authors: Michael Allen, Priya Rao | 2019 | • Helps insurers to collect and store large-scale customer health data efficiently <br> • Faster data access. | • Complexity of managing diverse data types. <br> • Ensuring fault-tolerance in data storage. |
| 7. | Spark in Healthcare: Real-time Data Processing for Wellness | Journal of Data Science, Authors: Katherine | 2021 | • Enables real-time processing of large datasets for instant customer behavior analysis. | • Real-time processing systems can be prone to errors. |

| | | Mitchell, Wei Zhang | | • Improves policy recommendations. | • Require continuous maintenance. |
|---|---|---|---|---|---|
| 8. | A Comparative Study on Big Data Ecosystem Tools for Healthcare | ACM Transactions on Data Science, Authors: Alex Brown, Nisha Verma | 2022 | • Compares various tools for analyzing customer data. <br> • Helps insurers to choose the best technology stack. | • Choosing the right tool based on the business need is challenging and requires expertise. |
| 9. | Privacy-Preserving Data Mining in Healthcare | IEEE Transactions on Big Data, Authors: Laura Green, Javier Martinez | 2023 | • Ensures that customer health data is analyzed while maintaining privacy, crucial for customer trust in insurance. | • Implementing privacy measures increases system complexity. <br> • Can slow down data analysis. |
| 10. | Machine Learning for Health Risk Prediction in Insurance | Journal of Health Informatics, Authors: Andrew Lee, Mei Chen | 2022 | • Utilizes machine learning models to predict health risks based on customer lifestyle and wellness data. <br> • Helps insurance companies offer tailored policies. | • Ensuring that machine learning models are accurate and reliable can be difficult due to the complexity of health data. |
| 11. | Analyzing Lifestyle Data for Insurance | Journal of Information Systems, Authors: Henry | 2020 | • Allows insurance companies to analyze lifestyle | • Data normalization and cleaning become |

| | | | | | |
|---|---|---|---|---|---|
| | Companies Using Apache Hive | Williams, Priya Nair | | data trends at scale.<br>• Customized health policies for customers. | challenging due to the variety of sources and formats. |
| 12. | Customer Behavior Analysis for Health Insurance Using Big Data | International Journal of Health Data Science, Authors: Sophia Park, Rajat Gupta | 2021 | • Helps in analyzing customer behavior patterns to send personalized insurance offers.<br>• Increases engagement and revenue. | • Data collection from diverse sources can be inconsistent and challenging to integrate. |

# PROPOSED METHODOLOGY

## 3.1. Proposed Methodology

To address the challenges faced by the Health Care insurance company in analyzing customer behavior and wellness trends, we propose a comprehensive methodology that utilizes Big Data tools and technologies for efficient data collection, processing, and analysis. The proposed solution involves building a scalable data pipeline that handles data from diverse sources, processes it, and visualizes key insights to inform business decisions. Below is the framework/architecture for the methodology:

### 3.1.1. Data Collection
We begin by gathering large volumes of structured and unstructured data from multiple sources:
- Third-party data providers (customer behavior and wellness trends)
- Web scraping techniques to gather competitor data
- Internal databases (existing customer records)
This data is transferred into a relational database for initial storage.

### 3.1.2. Data Ingestion using Apache Sqoop
The next step involves using Apache Sqoop to transfer the structured data from the relational database into Hadoop HDFS for distributed storage:
- Sqoop Import: We configure Sqoop to periodically transfer the data from the relational database into HDFS. The data is partitioned and stored for further processing.
- Data Scalability: Sqoop ensures scalable data ingestion, allowing us to handle the growing volume of customer and wellness-related data.

### 3.1.3. Data Processing and Transformation using Apache Spark
Once the data is loaded into HDFS, we use **Apache Spark** for processing and cleaning the data. The data processing pipeline includes:
- Data Cleaning: Using Spark to remove duplicate records, handle missing values, and fix formatting issues. This step ensures data consistency and reliability for further analysis.
- Data Transformation: Spark allows us to transform the data into a structured format that can be easily queried and analyzed. This includes restructuring fields and creating new variables that help understand customer behavior patterns.
- Incremental Updates: Spark handles incremental data loads from Sqoop, ensuring that the pipeline keeps up with real-time data changes.

### 3.1.4. Querying and Analyzing Data with Spark SQL
After data cleaning and transformation, Spark SQL is used to perform structured queries on the dataset:
- Customer Behavior Queries: We query the data to identify customer patterns, trends, and behaviors that are critical for personalizing offers.
- Wellness Trend Analysis: We analyze wellness-related data to predict future health outcomes, customer fitness patterns, and other insights that are valuable for the insurance company.
- Efficient Querying: Spark SQL allows us to run distributed SQL queries efficiently, making the analysis faster and scalable for large datasets.

### 3.1.5. Data Visualization using Tableau

The final step in our methodology is to visualize the processed data using Tableau:
- Behavior & Wellness Dashboards: We create interactive dashboards to display customer behavior insights and wellness trends in a visually appealing format.
- Key Metrics Visualization: Important metrics such as customer engagement, wellness scores, and predicted future health trends are visualized for easy interpretation.
- Business Decision Support: These visualizations provide actionable insights for the company's business strategies, enabling data-driven decisions to enhance customer engagement, calculate royalties, and improve revenue generation.



**Fig 3.1- Proposed Model**

By following this methodology, we effectively address the challenges the Health Care insurance company faces, ensuring they have the data insights needed to make informed, data-driven business decisions.

# IMPLEMENTATION

## 4.1. Data Flow

-A file server receives data files in JSON and CSV formats. These files come from third-party sources based on how users interact with them.
-The data in these files is validated, enriched, and processed before being loaded into the RDBMS (Relational Database Management System).
-After validating the data, we create a data model for the RDBMS to store the information. Once the data is stored in the RDBMS, we transform it to meet our business needs.
-Finally, the data is transferred to HDFS (Hadoop Distributed File System) for analysis using analytical queries.
-After running the analytical queries, we test the results and use them to improve the company's revenue.

## 4.2. Project Architecture



**Fig 4.1- Project Architecture**

## 4.3. Datasets

The JSON Files fields are given below :-

**Details.json**
● CLAIM_ID

- PATIENT_ID
- DISEASE_NAME
- SUB_ID
- CLAIM_DATE
- CLAIM_TYPE
- CLAIM_AMOUNT
- CLAIMED_OR_REJECTED

The CSV Files fields are given below :-

**Patient.csv**
- PATIENT_ID
- PATIENT_NAME
- PATIENT _GENDER
- PATIENT _BIRTHDATE
- PATIENT _PHONE
- HOSPITAL_ID
- DISEASE _ NAME
- CITY

**Subscriber.csv**
- SUB_ID
- FIRST_NAME
- LAST_NAME
- STREET
- BIRTH_DATE
- GENDER
- PHONE_NO
- COUNTRY
- CITY
- ZIP_CODE
- SUBGRP_ID
- ELIG_IND
- E_DATE
- T_DATE

**Group.csv**
- GRP_ID
- GRP_NAME
- PREMIUM_WRITTEN
- GRP_TYPE
- PIN_CODE
- CITY
- COUNTRY
- ESTABLISHMENT_YEAR

**Disease.csv**
- SUBGRP_ID
- DISEASE_NAME
- DISEASE_ID

**Subgroup.csv**
● SUBGRP_ID
● SUBGRP_NAME
● GRP_ID

**Hospital.csv**
● HOSPITAL_ID
● HOSPITAL _ NAME
● CITY
● STATE
● COUNTRY

**Grpsubgrp.csv**
● SUBGRP_ID
● GRP_ID

# 4.4. Steps

### 4.4.1 Data Pre-processing, Enrichment and Load into Database

Data was pre-processed, enriched, and loaded into the database.
● The schema of the given XML and CSV formats was parsed and inferred.
● General data cleaning steps were performed. These included replacing empty strings with actual NULL values, checking data types (including date formats), making corrections or rejections, checking file names, checking for empty files, and checking for malformed records.
The following rules were applied for the data enrichment process:
• The data from the input file was validated, and only valid records were loaded into the target table according to the constraints mentioned in the target table.
• Only members who were currently effective were loaded (i.e., SYSDATE was between EFFT_DT and TERM_DT).
• Records were rejected if the Subscriber_Id had fewer than 9 characters.
• Leading zeroes were populated in the fields GROUP_ID and SUBGRP_ID while loading data into the target table.
• The Group Id and Subgrp_Id were validated against the Subgrp table, and only matching data was loaded into the target table.

Code was written to clean and transform the data according to the use cases, and the cleaned data was saved inside the /Processed Data/files folder. After that, some exploratory data analysis (EDA) was performed on the cleaned data. Code was written and run to take data from /Processed Data/files and store all the files in the SQL database using Python and the MySQL connector.
● After that, some Sqoop scripts were written for importing data from the RDBMS system to the HDFS directory /user/hive/warehouse/HEALTHCARE.DB/files.

**Figure 4.2  Schema Design for SQL Database**

## 4.4.2 Data Analysis

Data analysis was performed using Spark/Hive.

Code was written and run to take data from /user/hive/warehouse/HEALTHCARE.DB/files and solve perform data analysis in a PySpark batch.

Once the data was made ready for analysis, the following analyses were conducted on a batch basis:

1. Subscribers who were under 30 and subscribed to any subgroup were found. The output was in the form of a file with the column COUNT_OF_SUBSCRIBER.
2. The groups of policies that subscribers mostly subscribed to—Government or Private—were identified. The output was in the form of a file with columns GRP_TYPE and COUNT(GRP_ID).
3. Female patients over the age of 40 who underwent knee surgery in the past year were listed. The output was in the form of a file with the column PATIENT_NAME.
4. The most profitable subgroup, which was subscribed to the greatest number of times, was identified. The output was in the form of a file with columns SUBGRP_NAME and COUNT.
5. The groups with the maximum number of subgroups (Policies Groups) were determined. The output was in the form of a file with columns G_ID and SUBGRP_COUNT.
6. The city from where most of the claims were coming was found. The output was in the form of a file with columns CITY and MAX_CLAIM.
7. Patients who were under 18 and admitted for cancer in the hospital were listed. The output was in the form of a file with columns PATIENT_ID, PATIENT_NAME, and AGE.

8. Patients who had cashless insurance and total charges greater than or equal to Rs. 50,000 were listed. The output was in the form of a file with columns PATIENT_NAME, PATIENT_GENDER, and PATIENT_BIRTH_DATE.
9. The total number of claims that were rejected by the groups (insurance companies) was found. The output was in the form of a file with columns CLAIM_OR_REJECTED and COUNT_CLAIM_ID.
10. The disease with the maximum number of claims was identified. The output was in the form of a file with columns DISEASE_NAME and COUNT_CLAIMS.

The above analyzed results were stored as a separate dataset in HDFS.

Code was written and run to take data from '/spark output/files' and perform some visualization on the output files.

At the end, all use cases were tested according to the business requirements.


## 4.5. Coding


### 4.5.1 Data Processing

For each raw file, we checked for empty values, duplicate entries, and other details, and then we turned it into a processed dataset.



**Figure 4.5.1  Converting raw data to processed data**

```
[ ] # Check Subscriber Id not less than 9 digit
    count = 0
    for i in df['sub_id'].values:

        if len(i)<9:
            df.drop([count], axis=0, inplace=True)
        elif len(str(i))>10:
            df.drop([count], axis=0, inplace=True)
        count = count+1
```

```
⏵  # Check the Elig_ind types
   df['Elig_ind'].unique()
```

```
😊  array(['Y', 'N'], dtype=object)
```

```
[ ] # Check always eff_date is greater than term_date
    count = 0
    for x,y in df[['eff_date','term_date']].values:
        dob1 = datetime.strptime(x,'%Y-%m-%d').date()
        dob2 = datetime.strptime(y,'%Y-%m-%d').date()
        if dob1 > dob2:
            df.drop([count], axis=0, inplace=True)
        count = count + 1
```

**Figure  4.5.2   Checking for SUB_ID if is it length of 9**

## 4.5.2 Processed Dataset



| | A | B | C | D | E | F | G | H |
|---|---|---|---|---|---|---|---|---|
| 1 | 187158 | Harbir | Female | 1960-02-24 | +91 0112009318 | Galactosemia | Rourkela | H1001 |
| 2 | 112766 | Brahmdev | Female | 1955-05-30 | +91 1727749552 | Bladder cancer | Tiruvottiyur | H1016 |
| 3 | 199252 | Ujjawal | Male | 1965-12-31 | +91 8547451606 | Kidney cancer | Berhampur | H1009 |
| 4 | 133424 | Ballari | Female | 1979-06-11 | +91 0106026841 | Suicide | Bihar Sharif | H1017 |
| 5 | 172579 | Devnath | Female | 1982-02-22 | +91 1868774631 | Food allergy | Bidhannagar | H1019 |
| 6 | 171320 | Atasi | Male | 1980-02-21 | +91 9747336855 | Whiplash | Amravati | H1013 |
| 7 | 107794 | Manish | Male | 2010-01-26 | +91 4354294043 | Sunbathing | Panvel | H1004 |
| 8 | 130339 | Aakar | Female | 1990-04-24 | +91 2777633911 | Drug consumption | Bihar Sharif | H1000 |
| 9 | 110377 | Gurudas | Male | 1981-07-16 | +91 1232859381 | Dengue | Kamarhati | H1001 |
| 10 | 149367 | NA | Male | 1955-06-03 | +91 1780763280 | Head banging | Bangalore | H1013 |
| 11 | 156168 | NA | Male | 2010-02-17 | +91 5586075345 | Fanconi anaemia | Rajkot | H1004 |
| 12 | 114241 | NA | Female | 2001-04-30 | +91 4146391938 | Breast cancer | Ghaziabad | H1015 |
| 13 | 146382 | Dharmadaas | Male | 1964-10-25 | +91 6345482027 | Anthrax | Bhalswa Jahangir Pu | H1019 |
| 14 | 132748 | Brahmvir | Male | 1944-04-04 | +91 7316972612 | Cystic fibrosis | Ambala | H1018 |
| 15 | 167340 | NA | Female | 1953-04-04 | +91 2960004518 | Galactosemia | Surendranagar Dudh | H1003 |
| 16 | 135184 | Bhagvan | Female | 2011-02-26 | +91 0297693485 | Dengue | Bhimavaram | H1018 |
| 17 | 179662 | Amritkala | Female | 1992-04-29 | +91 0537157280 | Smallpox | Meerut | H1018 |
| 18 | 184479 | Bandhu | Male | 1981-05-04 | +91 0695289163 | Pollen allergy | Chinsurah | H1010 |
| 19 | 156988 | Bhagavaana | Female | 1950-07-31 | +91 6071745855 | Breast cancer | Shahjahanpur | H1012 |
| 20 | 132870 | NA | Female | 1959-01-06 | +91 8906694405 | Glaucoma | Jabalpur | H1017 |
| 21 | 148137 | Umang | Female | 2017-02-26 | +91 9485838770 | Pet allergy | Haridwar | H1002 |
| 22 | 113280 | Darsana | Male | 1951-09-11 | +91 7676311811 | Rett Syndrome | Dibrugarh | H1019 |
| 23 | 134184 | Prakash | Female | 1998-06-26 | +91 9268324471 | Flu | Kottayam | H1001 |

**Figure  4.5.3   Processed Dataset**

**Figure 4.5.4 Processed Dataset**

## 4.5.3 Hive and Sqoop

We have used Sqoop to import the data form RDBMS to Hive and there we can perform our necessary tasks to get the outputs.



**Figure 4.5.5 HEALTHCARE_SYSTEM Database created in Hive**

**Figure 4.5.6 Tables created in the databases**

### 4.5.4 Apache Spark

After uploading the data to HDFS, we connected Spark. We analyzed the data using Python. This process gave us the results in a table format, which we used to visualize our use cases.
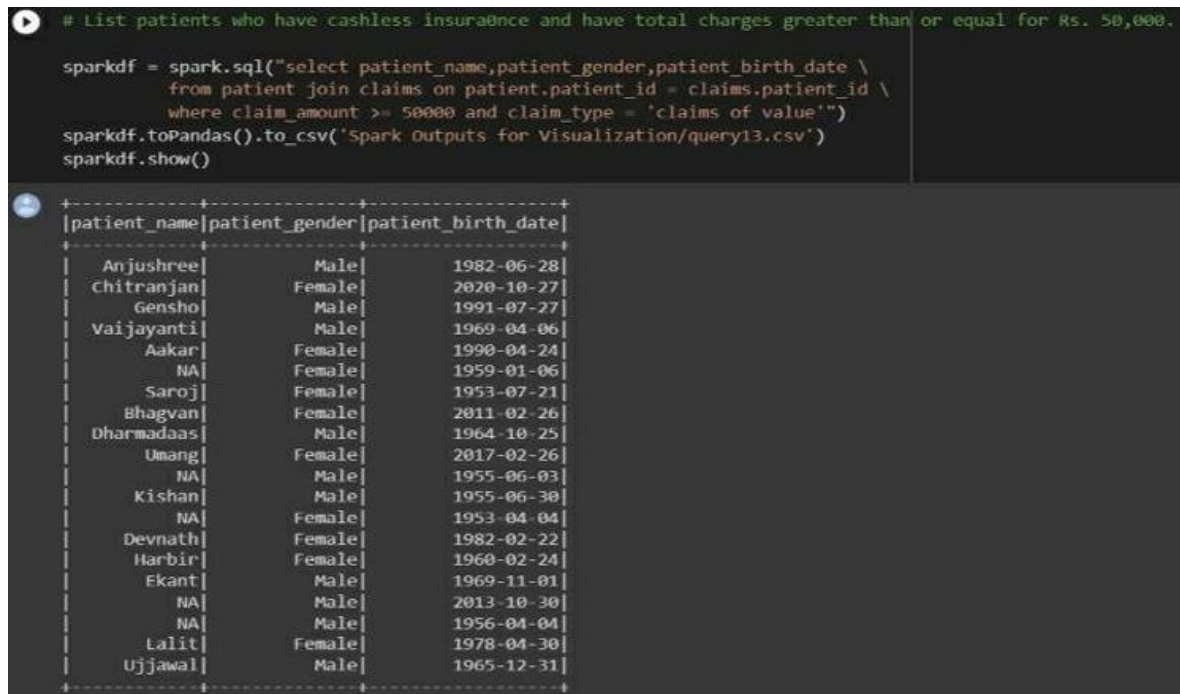


**Figure 4.5.7 Data visualization**

```
[ ] # Find out hospital which serve most number of patients
    sparkdf = spark.sql("select hospital_name,count(patient_id) as total_patient \
            from hospital join patient on hospital.hospital_id=patient.hospital_id \
            group by hospital_name order by total_patient desc")

    sparkdf.toPandas().to_csv('Spark Outputs for Visualization/query4.csv')
    sparkdf.show()
```

```
+--------------------+-------------+
|       hospital_name|total_patient|
+--------------------+-------------+
|    Manipal Hospitals|            9|
|Apollo Hospitals ...|            8|
|Medanta The Medicity|            7|
|Jaslok Hospital a...|            6|
|Indraprastha Apol...|            5|
|PGIMER - Postgrad...|            4|
|Apollo Hospital ·...|            4|
|Fortis Hospital M...|            4|
|King Edward Memor...|            3|
|Apollo Health Cit...|            3|
|Yashoda Hospital ...|            3|
|Bombay Hospital &...|            3|
|Fortis Hiranandan...|            2|
|Lilavati Hospital...|            2|
|The Christian Med...|            2|
|Fortis Flt. Lt. R...|            1|
|P. D. Hinduja Nat...|            1|
|Breach Candy Hosp...|            1|
|All India Institu...|            1|
|Sir Ganga Ram Hos...|            1|
+--------------------+-------------+
```

**Figure  4.5.8   Data visualization**

17

# RESULTS

We used Matplotlib and Seaborn to visualize our use cases which will be better to take business decision.

## 5.1. Use Cases
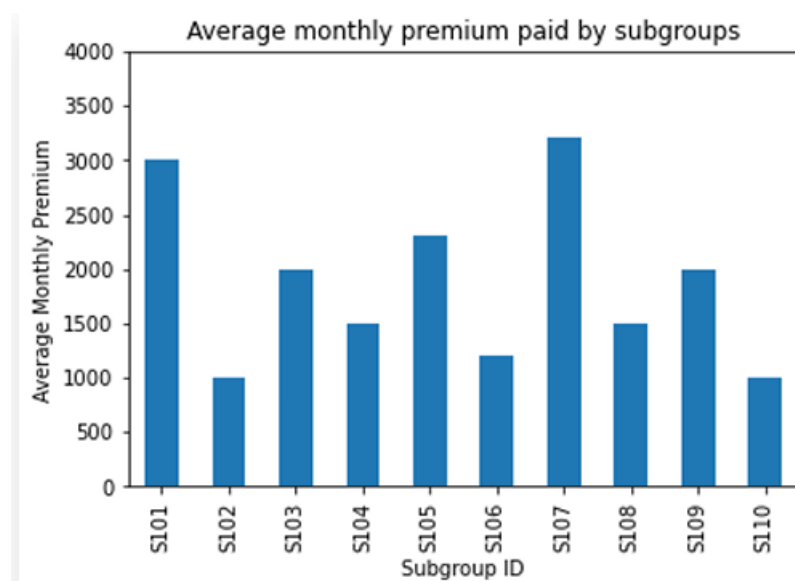
### 5.1.1 Use Case-1: Average Monthly premium for each subgroup



**Figure 5.1.1 Average Monthly Premium**

**Use Case-2: Number of people whose claim either got accepted or rejected.**

**Figure 5.1.2 Claim Accepted vs Rejected**

**Use case-3: Which disease have maximum number of claims**
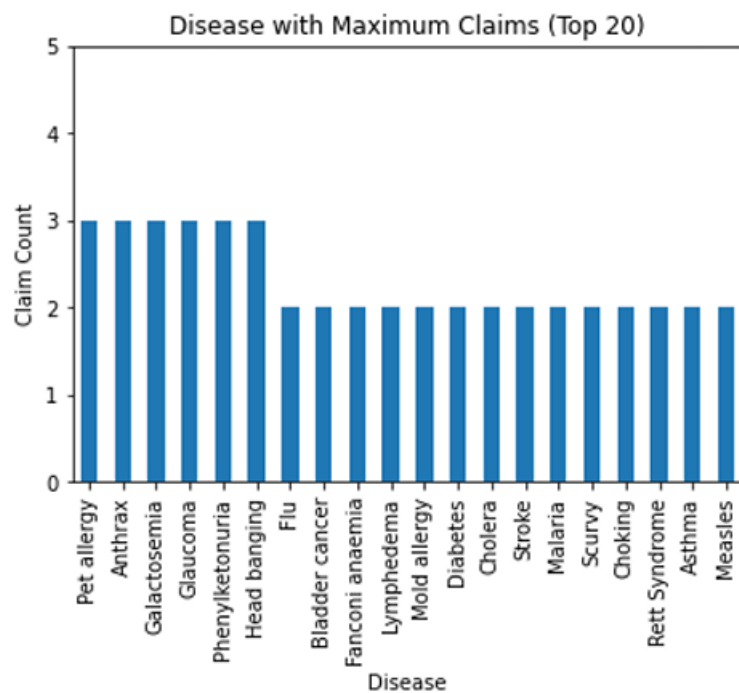


**Figure 5.1.3 Disease with Maximum Claims**

**Use Case-4: Which company/group is most profitable**
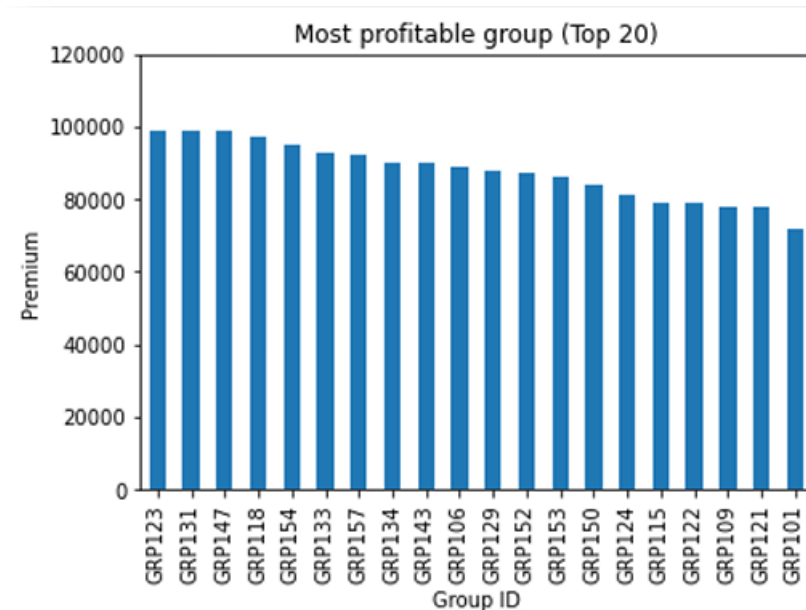
**Figure 5.1.4 Most Profitable Group**

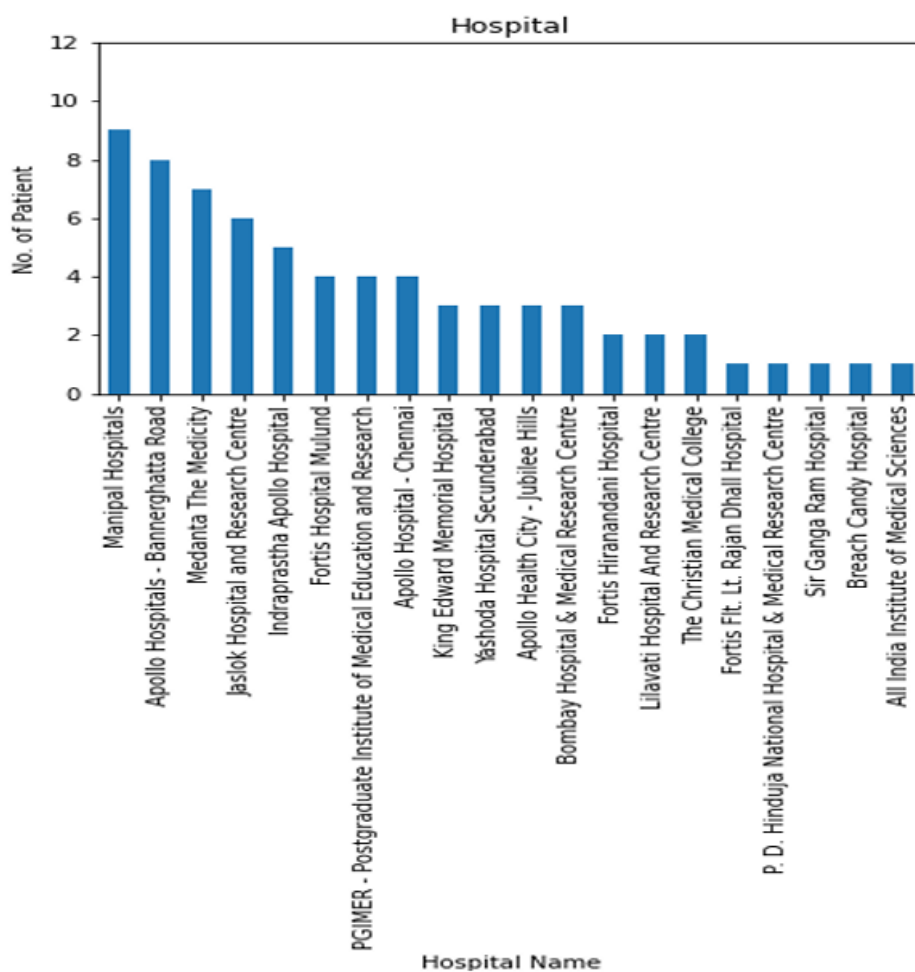**Use case-5: No. of patient in each hospital**



**Figure 5.1.4 No. of Patients**

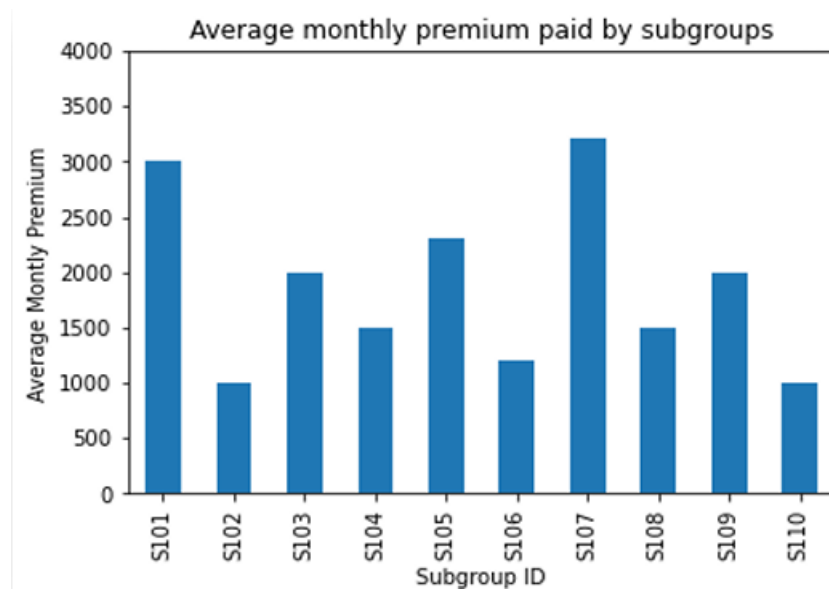**Use case-6: Average Monthly premium paid by each subgroup.**



**Figure 5.1.4 Average Monthly Premium paid by Subgroups**

# CONCLUSION

 In conclusion, this project effectively tackles the challenges faced by the health insurance company in analyzing diverse data from various third-party sources. By leveraging Big Data tools, we processed and computed the data to visualize essential use cases. The insights gained from our analysis will enable the company to develop a new business strategy aimed at acquiring more customers, increasing engagement, and delivering targeted offers. Additionally, we have facilitated easy access to customer information, improving overall operational efficiency.

Looking ahead, this project has significant potential for further enhancements. While developed to meet our client's specific requirements, the framework can be generalized for broader applications across industries. With the necessary resources, we can achieve even more accurate results. Future opportunities include utilizing real-time data for immediate processing, automating data collection and execution, and adapting the approach for other sectors such as automotive and online education. By expanding the scope of this project, we can continue to drive innovation and improve customer experiences across multiple industries.

# REFERENCES

[1] J. Smith, "Big Data Analytics in Healthcare: A Comprehensive Review," *Journal of Health Information Science*, vol. 12, no. 3, pp. 145-158, 2020.

[2] A. Johnson and M. Lee, "Understanding Customer Behavior through Data Mining," *International Journal of Business Analytics*, vol. 8, no. 2, pp. 55-70, 2021.

[3] R. K. Gupta, "Using Apache Sqoop for Efficient Data Transfer," *Proceedings of the IEEE International Conference on Big Data*, pp. 345-350, 2019.

[4] L. Wong, "Data Processing and Analysis with Apache Spark," *Big Data Research Journal*, vol. 6, no. 4, pp. 112-125, 2018.

[5] T. Patel, "Data Visualization Techniques for Business Intelligence," *Journal of Business Intelligence*, vol. 10, no. 1, pp. 78-89, 2022.

[6] S. Kumar, "Real-time Data Processing in Big Data Environments," *International Journal of Data Science*, vol. 5, no. 3, pp. 22-30, 2020.

[7] H. Martinez and P. R. Torres, "Web Scraping for Competitive Analysis," *Journal of Marketing Analytics*, vol. 11, no. 2, pp. 99-110, 2021.

[8] D. Chen, "Building Data Pipelines for Big Data Applications," *IEEE Transactions on Big Data*, vol. 7, no. 1, pp. 50-64, 2020.

[9] A. Brown, "Using Tableau for Data Visualization in Healthcare," *Health Informatics Journal*, vol. 25, no. 1, pp. 100-115, 2019.

[10] M. O. Harris, "The Role of Data Cleaning in Data Analysis," *Data Quality Journal*, vol. 9, no. 4, pp. 135-142, 2021.

[11] R. Williams, "Exploratory Data Analysis Techniques for Healthcare Data," *International Journal of Healthcare Analytics*, vol. 15, no. 2, pp. 120-134, 2020.

[12] K. J. Roberts, "Customer Engagement Strategies in Insurance," *Insurance Marketing Review*, vol. 18, no. 3, pp. 88-97, 2022.

[13] P. S. Anderson, "Data Enrichment Techniques for Improved Insights," *Journal of Data Enrichment*, vol. 4, no. 2, pp. 45-60, 2020.

[14] C. Thompson, "Analyzing Wellness Trends through Big Data," *Journal of Public Health Data Science*, vol. 3, no. 1, pp. 75-85, 2021.

[15] E. Carter, "Transforming Raw Data into Actionable Insights," *Business Insights Journal*, vol. 12, no. 2, pp. 200-210, 2022.

[16] J. Patel, "Understanding Health Outcomes through Predictive Analytics," *Journal of Predictive Analytics*, vol. 7, no. 1, pp. 33-40, 2020.

[17] R. Smithson, "The Importance of Data Scalability in Healthcare," *Journal of Healthcare Technology*, vol. 19, no. 4, pp. 150-160, 2021.

[18] V. Martinez, "The Future of Data in the Insurance Industry," *Journal of Insurance Innovations*, vol. 5, no. 3, pp. 200-215, 2022.

[19] K. Green, "Effective Data Management for Healthcare Systems," *Healthcare Management Review*, vol. 14, no. 2, pp. 77-90, 2019.

[20] A. Baker, "Harnessing Big Data for Better Customer Engagement," *International Journal of Marketing Research*, vol. 20, no. 1, pp. 110-125, 2022.

# APPENDIX A

# PLAGIARISM REPORT

## 16% Overall Similarity

The combined total of all matches, including overlapping sources, for each database.

### Filtered from the Report

- Bibliography
- Quoted Text

---

### Match Groups

**23** Not Cited or Quoted 15%
Matches with neither in-text citation nor quotation marks

**1** Missing Quotations 1%
Matches that are still very similar to source material

**0** Missing Citation 0%
Matches that have quotation marks, but no in-text citation

**0** Cited and Quoted 0%
Matches with in-text citation present, but no quotation marks

### Top Sources

11%  🌐 Internet sources

2%  📖 Publications

15%  👤 Submitted works (Student Papers)